

# Embedding Methods for NLP

## Part 1: Unsupervised and Supervised Embeddings

Jason Weston & Antoine Bordes

Facebook, NY, USA

# What is a word embedding?

Suppose you have a dictionary of words.

The  $i^{th}$  word in the dictionary is represented by an embedding:

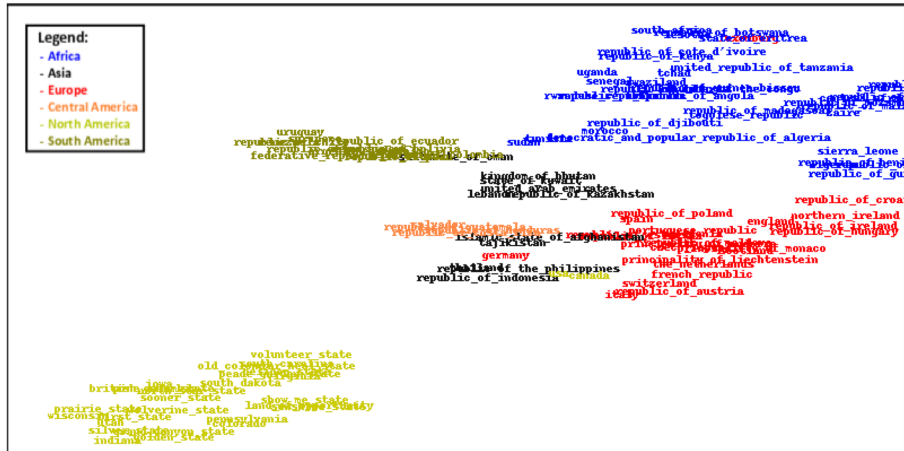
$$w_i \in \mathbb{R}^d$$

i.e. a  $d$ -dimensional vector, which is **learnt!**

- $d$  typically in the range 50 to 1000.
- Similar words should have similar embeddings (share latent features).
- Embeddings can also be applied to *symbols* as well as words (e.g. Freebase nodes and edges).
- Discuss later: can also have embeddings of phrases, sentences, documents, or even other modalities such as images.

# Learning an Embedding Space

## Example of Embedding of 115 Countries (Bordes et al., '11)



## Main methods we highlight, ordered by date.

- Latent Semantic Indexing (Deerwester et al., '88).
- Neural Net Language Models (NN-LMs) (Bengio et al., '06)
- Convolutional Nets for tagging (SENNA) (Collobert & Weston, '08).
- Supervised Semantic Indexing (Bai et al, '09).
- Wsabee (Weston et al., '10).
- Recurrent NN-LMs (Mikolov et al., '10).
- Recursive NNs (Socher et al., '11).
- Word2Vec (Mikolov et al., '13).
- Paragraph Vector (Le & Mikolov, '14).
- Overview of recent applications.

# Main methods we highlight, ordered by topic.

## Embeddings for Ranking and Retrieval:

- Latent Semantic Indexing (Deerwester et al., '88).
- Supervised Semantic Indexing (Bai et al, '09).
- Wsabee (Weston et al., '10).

## Embeddings for Language Modeling (useful for speech, translation, ...):

- Neural Net Language Models (NN-LMs) (Bengio et al., '06)
- Recurrent NN-LMs (Mikolov et al., '10).
- Word2Vec (Mikolov et al., '13).

## Embeddings for Supervised Prediction Tasks (POS, chunk, NER, SRL, sentiment, etc.):

- Convolutional Nets for tagging (SENNA) (Collobert & Weston, '08).
- Recursive NNs (Socher et al., '11).
- Paragraph Vector (Le & Mikolov, '14).

# Main methods we highlight, ordered by topic.

## Embeddings for Ranking and Retrieval:

- Latent Semantic Indexing (Deerwester et al., '88).
- Supervised Semantic Indexing (Bai et al, '09).
- Wsabee (Weston et al., '10).

## Embeddings for Language Modeling (useful for speech, translation, ...):

- Neural Net Language Models (NN-LMs) (Bengio et al., '06)
- Recurrent NN-LMs (Mikolov et al., '10).
- Word2Vec (Mikolov et al., '13).

## Embeddings for Supervised Prediction Tasks (POS, chunk, NER, SRL, sentiment, etc.):

- Convolutional Nets for tagging (SENNA) (Collobert & Weston, '08).
- Recursive NNs (Socher et al., '11).
- Paragraph Vector (Le & Mikolov, '14).

# Ranking and Retrieval: The Goal

We want to learn to match a query (text) to a target (text).



Many classical supervised ranking methods use hand-coded features.



Methods like LSI that learn from words are unsupervised.



Supervised Semantic Indexing (SSI) uses supervised learning from text only:

*Bai et al, Learning to Rank with (a Lot of) Word Features. Journal of Information Retrieval, '09.*

*Outperforms existing methods (on words) like TFIDF, LSI or a (supervised) margin ranking perceptron baseline.*

## Basic Bag-O'-Words



Bag-of-words + cosine similarity:

- Each doc.  $\{d_t\}_{t=1}^N \subset \mathbb{R}^{\mathcal{D}}$  is a *normalized* bag-of-words.
- Similarity with query  $q$  is:  $f(q, d) = q^\top d$



**Doesn't deal with synonyms:** bag vectors can be orthogonal



**No machine learning at all**



# Latent semantic indexing (LSI)



Learn a linear embedding  $\phi(d_i) = Ud_i$  via a reconstruction objective.

- Rank with:  $f(q, d) = q^\top U^\top U d = \phi(q)^\top \phi(d)$ <sup>1</sup>.



Uses “synonyms”: *low-dimensional latent “concepts”*.



Unsupervised machine learning: useful for goal?

---

<sup>1</sup> $f(q, d) = q^\top (U^\top U + \alpha I) d$  gives better results.  
Also, usually normalize this  $\rightarrow$  cosine similarity.

# Supervised Semantic Indexing (SSI)

- Basic model: rank with

$$f(q, d) = q^T W d = \sum_{i,j=1}^D q_i W_{ij} d_j$$

i.e. learn weights of polynomial terms between documents.

- Learn  $W \in \mathbb{R}^{D \times D}$  (huge!) with click-through data or other labels.



Uses “synonyms”



Supervised machine learning: targeted for goal



Too Big/Slow?! Solution = Constrain  $W$  :

**low rank**  $\rightarrow$  embedding model!

# SSI: why is this a good model?

Classical bag-of-words doesnt work when there are few matching terms:

$q = (\text{kitten}, \text{vet}, \text{nyc})$

$d = (\text{cat}, \text{veterinarian}, \text{new}, \text{york})$



Method  $q^T W d$  learns that e.g. **kitten** and **cat** are highly related.



E.g. if  $i$  is the index of **kitten** and  $j$  is the index of **cat**, then  $W_{ij} > 0$  after training.

# SSI: Why the Basic Model Sucks



$W$  is **big** : 3.4Gb if  $\mathcal{D} = 30000$ , 14.5Tb if  $\mathcal{D} = 2.5M$ .



**Slow**:  $q^\top W d$  computation has  $mn$  computations  $q_j W_{ij} d_i$ , where  $q$  and  $d$  have  $m$  and  $n$  nonzero terms.



Or one computes  $v = q^\top W$  once, and then  $vd$  for each document. Classical speed where query has  $\mathcal{D}$  terms, assuming  $W$  is dense  $\rightarrow$  **still slow**.



One could minimize  $\|W\|_1$  and attempt to make  $W$  sparse. Then at most  $mp$  times slower than classical model (with  $p$  nonzeros in a column.)

# SSI Improved model: Low Rank $W$

🗨️✎️ **Constrain  $W$ :**

$$W = U^T V + I.$$

🗨️✎️  $U$  and  $V$  are  $N \times \mathcal{D}$  matrices  $\rightarrow$  **smaller**

🗨️✎️ Low dimensional “latent concept” space like LSI (same speed).

🗨️✎️ Differences: supervised, asymmetric, learns with  $I$ .

Variants:

- $W = I$ : **bag-of-words** again.
- $W = D$ , **reweighted bag-of-words** related to [Grangier and Bengio, 2005].
- $W = U^T U + I$ : **symmetric**.

## SSI: Training via maximizing AUC

- Given a set of tuples  $\mathcal{R}$  with a query  $q$ , a related document  $d^+$  and an unrelated (or lower ranked) document  $d^-$ .
- We would like  $f(q, d^+) > f(q, d^-)$ .
- Minimize margin ranking loss [Herbrich et al., 2000]:

$$\sum_{(q, d^+, d^-) \in \mathcal{R}} \max(0, 1 - f(q, d^+) + f(q, d^-)).$$

Learning Algorithm Stochastic Gradient Descent: **Fast & scalable.**

Iterate	Sample a triplet $(q, d^+, d^-)$ , Update $W \leftarrow W - \lambda \frac{\partial}{\partial W} \max(0, 1 - f(q, d^+) + f(q, d^-))$ .
---------	--

Other options: batch gradient, parallel SGD (hogwild), Adagrad ...

## Prior Work: Summary of learning to Rank

- [Grangier & Bengio, '06] used similar methods to basic SSI for retrieving images.
- [Goel, Langford & Strehl, '08] used Hash Kernels (Vowpal Wabbit) for advert placement.
- Main difference: we use low rank+CFH on  $W$ .
- SVM [Joachims, 2002] and NN ranking methods [Burges, 2005] .  
Use hand-coded features: title, body, URL, search rankings,... (don't use words)  
(e.g. Burges uses 569 features in all).
- In contrast we use only the words and train on huge feature sets.
- Several works on optimizing different loss functions (MAP, ROC, NDCG):  
[Cao, 2008], [Yu, 2007], [Qin, 2006],...
- Lots of stuff for “metric learning” problem as well..



**One could also add features + new loss to this method ..**

# Experimental Comparison

- **Wikipedia**

- 1,828,645 documents. 24,667,286 links.
- Split into 70% train, 30% test.
- Pick random doc. as query, then rank other docs.
- Docs that are linked to it should be highly ranked.
- **Two setups:**
  - (i) whole document is used as query;
  - (ii) 5,10 or 20 words are picked to mimic keyword search.



# Wikipedia Experiments: Document Retrieval Performance

Experiments on Wikipedia, which still contains 2 million documents: retrieval task using the link structure and separated the data into 70% for training and 30% for test.

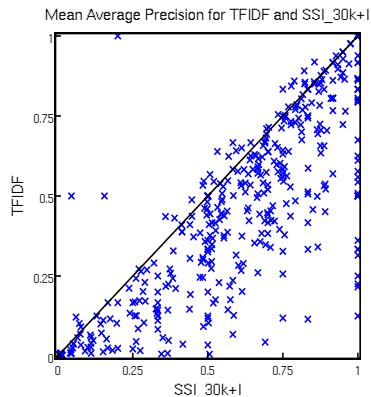
## Document based retrieval:

Algorithm	Rank-Loss	MAP	P10
TFIDF	0.842%	0.432 $\pm$ 0.012	0.193
$\alpha$ LSI + $(1 - \alpha)$ TFIDF	0.721%	0.433	0.193
Linear SVM Ranker	0.410%	0.477	0.212
Hash Kernels + $\alpha$ <i>l</i>	0.322%	0.492	0.215
Wsabie + $\alpha$ <i>l</i>	0.158%	0.547 $\pm$ 0.012	0.239 $\pm$ 0.008

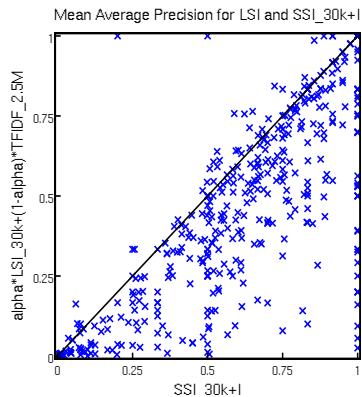
## *k*-keywords based retrieval:

<i>k</i> = 5: Algorithm	Params	Rank	MAP	P@10
TFIDF	0	21.6%	0.047	0.023
$\alpha$ LSI + $(1 - \alpha)$ TFIDF	200 $\mathcal{D}$ +1	14.2%	0.049	0.023
Wsabie + $\alpha$ <i>l</i>	400 $\mathcal{D}$	<b>4.37%</b>	<b>0.166</b>	<b>0.083</b>

# Scatter Plots: SSI vs. TFIDF and LSI



(a)



(b)

**Figure :** Scatter plots of Average Precision for 500 documents:  
(a) SSI vs. TFIDF, (b) SSI vs.  $\alpha$ LSI +  $(1 - \alpha)$  TFIDF.

## Experiments: Cross-Language Retrieval

Retrieval experiments using a query document in japanese, where the task is to retrieve documents in English (using link structure as ground truth).

SSI can do this without doing a translation step first as it learns to map the two languages together in the embedding space.

$\mathcal{D} = 30,000$

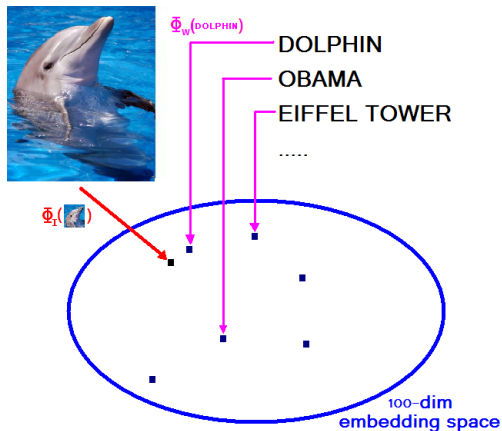
Algorithm	Rank-Loss	MAP	P10
TFIDF <sub>EngEng</sub> (Google translated queries)	4.78%	0.319	0.259
$\alpha$ LSI <sub>EngEng</sub> + $(1 - \alpha)$ TFIDF <sub>EngEng</sub>	3.71%	0.300	0.253
$\alpha$ CL-LSI <sub>JapEng</sub> + $(1 - \alpha)$ TFIDF <sub>EngEng</sub>	3.31%	0.275	0.212
SSI <sub>EngEng</sub>	1.72%	0.399	0.325
SSI <sub>JapEng</sub>	0.96%	0.438	0.351
$\alpha$ SSI <sub>JapEng</sub> + $(1 - \alpha)$ TFIDF <sub>EngEng</sub>	0.75%	0.493	0.377
$\alpha$ SSI <sub>JapEng</sub> + $(1 - \alpha)$ SSI <sub>EngEng</sub>	<b>0.63%</b>	<b>0.524</b>	<b>0.386</b>

Some recent related translation-based embeddings:

(Hermann & Blunsom, ICLR '14) and (Mikolov et al., '13).

# Wsabie (Weston, Usunier, Bengio, '10)

- Extension to SSI, also embeds objects other than text, e.g. images.
- WARP loss function that optimizes precision@k.



*Learn  $\Phi_I(\cdot)$  and  $\Phi_W(\cdot)$  to optimize precision@k.*

## Joint Item-Item Embedding Model

L.H.S: Image, query string or user profile (depending on the task)

$$\Phi_{LHS}(x) = U\Phi_x(x) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{100}.$$

R.H.S: document, image, video or annotation (depending on the task)

$$\Phi_{RHS}(y) = V\Phi_y(y) : \mathbb{R}^{d_y} \rightarrow \mathbb{R}^{100}.$$

This model again compares the degree of match between the L.H.S and R.H.S in the embedding space:

$$f_y(x) = \text{sim}(\Phi_{LHS}(x), \Phi_{RHS}(y)) = \Phi_x(x)^\top U^\top V\Phi_y(y)$$

*Also constrain the weights (regularize):*

$$\|U_i\|_2 \leq C, \quad i = 1, \dots, d_x, \quad \|V_i\|_2 \leq C, \quad i = 1, \dots, d_y.$$

# Ranking Annotations: AUC is Suboptimal

Classical approach to learning to rank is maximize AUC by minimizing:

$$\sum_x \sum_y \sum_{\tilde{y} \neq y} \max(0, 1 + f_{\tilde{y}}(x) - f_y(x))$$

**Problem:** All pairwise errors are considered the same, it counts the number of ranking violations.

**Example:**

Function 1: true annotations ranked 1st and 101st.

Function 2: true annotations ranked 50th and 52nd.

AUC prefers these *equally* as both have 100 “violations”.

**We want to optimize the top of the ranked list!**

# Rank Weighted Loss [Usunier et al. '09]

Replace classical AUC optimization:

$$\sum_x \sum_y \sum_{\bar{y} \neq y} \max(0, 1 + f_{\bar{y}}(x) - f_y(x))$$

With weighted version:

$$\sum_x \sum_y \sum_{\bar{y} \neq y} L(\text{rank}_y(x)) \max(0, 1 + f_{\bar{y}}(x) - f_y(x))$$

where  $\text{rank}_y(f(x))$  is the rank of the true label:

$$\text{rank}_y(f(x)) = \sum_{\bar{y} \neq y} I(f_{\bar{y}}(x) \geq f_y(x))$$

and  $L(\eta)$  converts the rank to a weight, e.g.  $L(\eta) = \sum_{i=1}^{\eta} 1/i$ .

# Weighted Approximate-Rank Pairwise (WARP) Loss

**Problem:** we would like to apply SGD:

$$\text{Weighting } L(\text{rank}_y(f(x))), \quad \text{rank}_y(f(x)) = \sum_{\bar{y} \neq y} I(f_{\bar{y}}(x) + 1 \geq f_y(x))$$

... too expensive to compute per  $(x, y)$  sample as  $y \in \mathcal{Y}$  is large.

**Solution:** approximate by sampling  $f_i(x)$  until we find a violating label  $\bar{y}$

$$\text{rank}_y(f(x)) \approx \left\lfloor \frac{|\mathcal{Y}| - 1}{N} \right\rfloor$$

where  $N$  is the number of trials in the sampling step.



# Online WARP Loss

**Input:** labeled data  $(x_i, y_i)$ ,  $y_i \in \{1, \dots, Y\}$ .

**repeat**

Pick a random labeled example  $(x_i, y_i)$

Set  $N = 0$ .

**repeat**

Pick a random annotation  $\bar{y} \in \{1, \dots, Y\} \setminus y_i$ .

$N = N + 1$ .

**until**  $f_{\bar{y}}(x) > f_{y_i}(x) - 1$  or  $N > Y - 1$

**if**  $f_{\bar{y}}(x) > f_{y_i}(x) - 1$  **then**

Make a **gradient step** to minimize:

$$L(\lfloor \frac{Y-1}{N} \rfloor) |1 - f_y(x) + f_{\bar{y}}(x)|_+$$

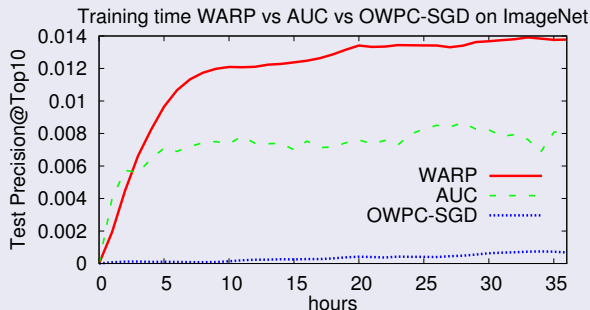
**end if**

**until** validation error does not improve.

## Image Annotation Performance

Algorithm	16k ImageNet	22k ImageNet	97k Web Data
Nearest Means	4.4%	2.7%	2.3%
One-vs-all SVMs 1+:1-	4.1%	3.5%	1.6%
One-vs-all SVMs	9.4%	8.2%	6.8%
AUC SVM Ranker	4.7%	5.1%	3.1%
Wsabie	11.9%	10.5%	8.3%

## Training time: WARP vs. OWPC-SGD & AUC



# Learned Annotation Embedding (on Web Data)

Annotation	Neighboring Annotations
barack obama david beckham santa	<i>barak obama, obama, barack, barrack obama, bow wow beckham, david beckam, alessandro del piero, del piero santa claus, papa noel, pere noel, santa clause, joyeux noel</i>
dolphin cows	delphin, dauphin, <i>whale, delfin, delfini, baleine, blue whale cattle, shire, dairy cows, kuh, horse, cow, shire horse, kone</i>
rose pine tree	rosen, <i>hibiscus, rose flower, rosa, roze, pink rose, red rose abies alba, abies, araucaria, pine, neem tree, oak tree</i>
mount fuji eiffel tower	mt fuji, fuji, fujisan, fujiyama, <i>mountain, zugspitze eiffel, tour eiffel, la tour eiffel, big ben, paris, blue mosque</i>
ipod f18	i pod, <i>ipod nano, apple ipod, ipod apple, new ipod f 18, eurofighter, f14, fighter jet, tomcat, mig 21, f 16</i>

# Summary

## Conclusion



**Powerful:** supervised methods for ranking.

- **Outperform** classical methods



**Efficient** low-rank models → learn hidden representations.



Embeddings good for generalization, but can “blur” too much e.g. for exact word matches.

## Extensions

- Nonlinear extensions – e.g. convolutional net instead.

# Main methods we highlight, ordered by topic.

## Embeddings for Ranking and Retrieval:

- Latent Semantic Indexing (Deerwester et al., '88).
- Supervised Semantic Indexing (Bai et al, '09).
- Wsabee (Weston et al., '10).

## Embeddings for Language Modeling (useful for speech, translation)

- Neural Net Language Models (NN-LMs) (Bengio et al., '06)
- Recurrent NN-LMs (Mikolov et al., '10).
- Word2Vec (Mikolov et al., '13).

## Embeddings for Supervised Prediction Tasks (POS, chunk, NER, SRL, sentiment, etc.):

- Convolutional Nets for tagging (SENNA) (Collobert & Weston, '08).
- Recursive NNs (Socher et al., '11).
- Paragraph Vector (Le & Mikolov, '14).

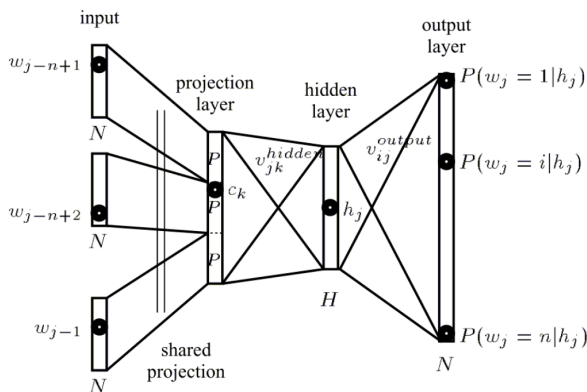
# Language Modeling

Task: given a sequence of words, predict the next word.

the cat sat on the ??

- $n$ -gram models are a strong baseline on this task.
- A variety of embedding models have been tried, they can improve results.
- The embeddings learnt from this unsupervised task can also be used to transfer to and improve a supervised task.

# Neural Network Language Models



Bengio, Y., Schwenk, H., Sencal, J. S., Morin, F., & Gauvain, J. L. (2006). Neural probabilistic language models. In *Innovations in Machine Learning* (pp. 137-186). Springer Berlin Heidelberg.

# Neural Network Language Models:

## Hierarchical Soft Max Trick (Morin & Bengio '05)

Predicting the probability of each next word is slow in NNLMs because the output layer of the network is the size of the dictionary.

Can predict via a tree instead:

- 1 Cluster the dictionary either according to semantics (similar words in the same cluster) or frequency (common words in the same cluster).  
*This gives a two-layer tree, but a binary tree is another possibility.*
- 2 The internal nodes explicitly model the probability of its child nodes.
- 3 The cost of predicting the probability of the true word is now: traversal to the child, plus normalization via the internal nodes and children in the same node.

*This idea is used in Word2Vec and RNN models as well.*



# Recurrent Neural Network Language Models

**Key idea:** *input to predict next word is current word plus context fed-back from previous word (i.e. remembers the past with recurrent connection).*

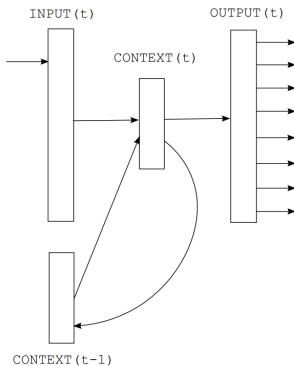


Figure: Recurrent neural network based LM

Recurrent neural network based language model. Mikolov et al., Interspeech, '10.

# NNLMS vs. RNNs: Penn Treebank Results (Mikolov)

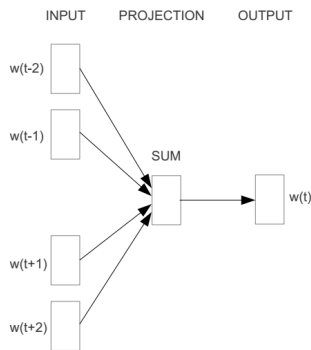
Model	Weight	PPL
3-gram with Good-Turing smoothing (GT3)	0	165.2
5-gram with Kneser-Ney smoothing (KN5)	0	141.2
5-gram with Kneser-Ney smoothing + cache	0.0792	125.7
Maximum entropy model	0	142.1
Random clusterings LM	0	170.1
Random forest LM	0.1057	131.9
Structured LM	0.0196	146.1
Within and across sentence boundary LM	0.0838	116.6
Log-bilinear LM	0	144.5
Feedforward NNLM	0	140.2
Syntactical NNLM	0.0828	131.3
Combination of static RNNLMs	0.3231	102.1
Combination of adaptive RNNLMs	0.3058	101.0
ALL	1	<b>83.5</b>

*Recent uses of NNLMs and RNNs to improve machine translation:*

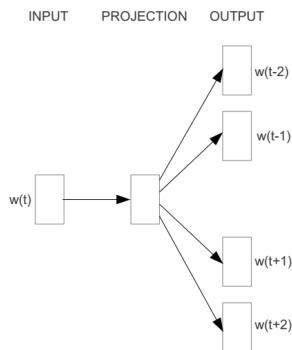
Fast and Robust NN Joint Models for Machine Translation, Devlin et al, ACL '14.

Also Kalchbrenner '13, Sutskever et al., '14., Cho et al., '14. .

# Word2Vec : very simple LM, works well



**CBOW**



**Skip-gram**

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean.  
Distributed Representations of Words and Phrases and their Compositionality.  
NIPS, 2013.

# Word2Vec: compositionality

Table 7: *Comparison and combination of models on the Microsoft Sentence Completion Challenge.*

Architecture	Accuracy [%]
4-gram [32]	39
Average LSA similarity [32]	49
Log-bilinear model [24]	54.8
RNNLMs [19]	55.4
Skip-gram	48.0
Skip-gram + RNNLMs	<b>58.9</b>

Czech + currency	Vietnam + capital	German + airlines	Russian + river	French + actress
koruna	Hanoi	airline Lufthansa	Moscow	Juliette Binoche
Check crown	Ho Chi Minh City	carrier Lufthansa	Volga River	Vanessa Paradis
Polish zolty	Viet Nam	flag carrier Lufthansa	upriver	Charlotte Gainsbourg
CTK	Vietnamese	Lufthansa	Russia	Cecile De

Table 5: Vector compositionality using element-wise addition. Four closest tokens to the sum of two vectors are shown, using the best Skip-gram model.

Code: <https://code.google.com/p/word2vec/>

# Main methods we highlight, ordered by topic.

## Embeddings for Ranking and Retrieval:

- Latent Semantic Indexing (Deerwester et al., '88).
- Supervised Semantic Indexing (Bai et al, '09).
- Wsabee (Weston et al., '10).

## Embeddings for Language Modeling (useful for speech, translation, ...):

- Neural Net Language Models (NN-LMs) (Bengio et al., '06)
- Recurrent NN-LMs (Mikolov et al., '10).
- Word2Vec (Mikolov et al., '13).

## Embeddings for Supervised Prediction Tasks (POS, chunk, NER, SRL, sentiment, etc.):

- Convolutional Nets for tagging (SENNA) (Collobert & Weston, '08).
- Recursive NNs (Socher et al., '11).
- Paragraph Vector (Le & Mikolov, '14).

# NLP Tasks

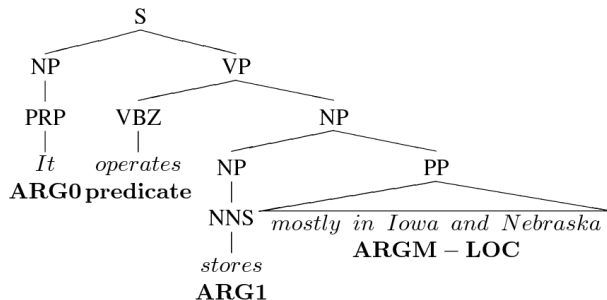
- Part-Of-Speech Tagging (POS): syntactic roles (noun, adverb...)
- Chunking: syntactic constituents (noun phrase, verb phrase...)
- Name Entity Recognition (NER): person/company/location...
- Semantic Role Labeling (SRL):

[John]*ARG0* [ate]*REL* [the apple]*ARG1* [in the garden]*ARGM-LOC*

## Complex Systems

- Two extreme choices to get a *complex system*
  - ★ *Large Scale Engineering*: *design* a lot of *complex features*, use a fast existing linear machine learning algorithm
  - ★ *Large Scale Machine Learning*: use simple features, design a *complex model* which will *implicitly learn* the right features

# The Large Scale Feature Engineering Way



- Extract **hand-made features** e.g. from the parse tree
- Disjoint: all tasks trained separately, Cascade features
- Feed these features to a shallow classifier like SVM

# ASSERT: many hand built features for SRL (Pradhan et al, '04)

## Problems:

- 1) Features rely on other solutions (parsing, named entity, word-sense)
  - 2) Technology task-transfer is **difficult**
- Choose some **good hand-crafted features**

**Predicate and POS tag** of predicate

**Phrase type**: adverbial phrase, prepositional phrase, ...

**Head word** and POS tag of the head word

**Path**: traversal from predicate to constituent

**Word-sense** disambiguation of the verb

**Length** of the target constituent (number of words)

**Partial Path**: lowest common ancestor in path

**First and last words** and POS in constituents

**Constituent tree distance**

**Dynamic class context**: previous node labels

**Constituent relative features**: head word

**Constituent relative features**: siblings

**Voice**: active or passive (hand-built rules)

**Governing category**: Parent node's phrase type(s)

**Position**: left or right of verb

Predicted **named entity** class

**Verb clustering**

**NEG feature**: whether the verb chunk has a "not"

**Head word replacement** in prepositional phrases

**Ordinal position** from predicate + constituent type

**Temporal cue words** (hand-built rules)

**Constituent relative features**: phrase type

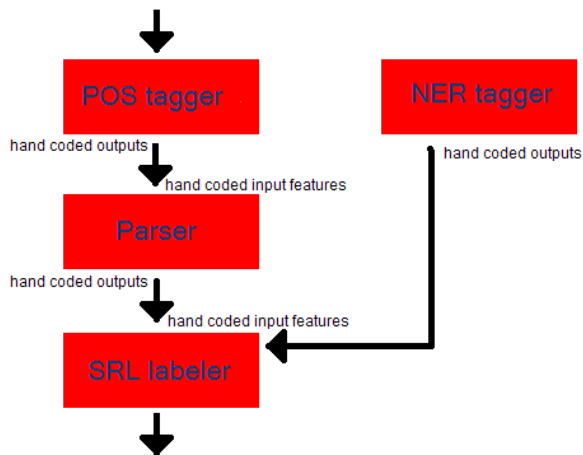
**Constituent relative features**: head word POS

Number of pirates existing in the world...

- Feed them to a **shallow classifier** like SVM



# The Suboptimal (?) Cascade



# NLP: Large Scale Machine Learning

## Goals

- Task-specific engineering **limits NLP scope**
- Can we find **unified hidden representations**?
- Can we build **unified NLP architecture**?

## Means

- Start **from scratch**: forget (most of) NLP knowledge
- Compare against classical **NLP benchmarks**
- **Our dogma**: avoid task-specific engineering

## NLP Benchmarks

- Datasets:

- ★ POS, CHUNK, SRL: [WSJ](#) ( $\approx$  up to 1M labeled words)
- ★ NER: [Reuters](#) ( $\approx$  200K labeled words)

System	Accuracy
Shen, 2007	97.33%
<b>Toutanova, 2003</b>	<b>97.24%</b>
Gimenez, 2004	97.16%

(a) **POS**: As in (Toutanova, 2003)

System	F1
<b>Ando, 2005</b>	<b>89.31%</b>
<a href="#">Florian</a> , 2003	88.76%
Kudoh, 2001	88.31%

(c) **NER**: CoNLL 2003

System	F1
Shen, 2005	95.23%
<b>Sha, 2003</b>	<b>94.29%</b>
Kudoh, 2001	93.91%

(b) **CHUNK**: CoNLL 2000

System	F1
<b>Koomen, 2005</b>	<b>77.92%</b>
Pradhan, 2005	77.30%
Haghighi, 2005	77.04%

(d) **SRL**: CoNLL 2005

- We chose as [benchmark systems](#):

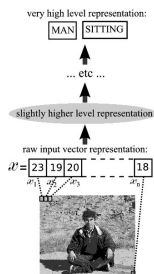
- ★ [Well-established](#) systems
- ★ Systems avoiding [external labeled](#) data

- Notes:

- ★ [Ando, 2005](#) uses external [unlabeled](#) data
- ★ [Koomen, 2005](#) uses 4 parse trees not provided by the challenge

# The “Deep Learning” Way

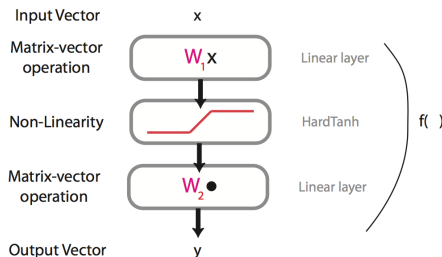
Deep approach attempts to propose a radically different **end-to-end** approach:



- **Avoid** building a **parse tree**. Humans don't need this to talk.
- Try to **avoid** all **hand-built** features → **monolithic systems**.
- Humans **implicitly** learn these features. Neural networks can too...?

## Neural Networks

- Stack several layers together



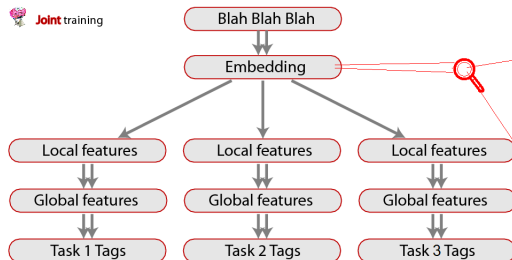
- Increasing level of abstraction at each layer
- Requires simpler features than “shallow” classifiers
- The “weights”  $W_i$  are trained by gradient descent
- How can we feed words?

# The Big Picture

A unified architecture for all NLP (labeling) tasks:

Sentence:	<i>Felix</i>	<i>sat</i>	<i>on</i>	<i>the</i>	<i>mat</i>	.
POS:	NNP	VBD	IN	DT	NN	.
CHUNK:	NP	VP	PP	NP	NP-I	.
NER:	PER	-	-	-	-	-
SRL:	ARG1	REL	ARG2	ARG2-I	ARG2-I	-

 Joint training



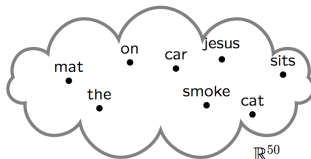
- ★ The *cat* sat on the mat
- ★ The *feline* sat on the mat



## Words into Vectors

### Idea

- Words are **embed** in a **vector** space



- Embeddings are **trained**

### Implementation

- A word  $w$  is an **index** in a dictionary  $\mathcal{D} \in \mathbb{N}$
- Use a **lookup-table** ( $W \sim \text{feature size} \times \text{dictionary size}$ )

$$LT_W(w) = W_{\bullet w}$$

### Remarks

- Applicable to any **discrete feature** (words, caps, stems...)
- See (Bengio et al, 2001)

# The Lookup Tables

Each word/element in dictionary maps to a vector in  $\mathbb{R}^d$ .

- We learn these vectors.
- LookupTable: input of  $i^{th}$  word is

$$x = (0, 0, \dots, 1, 0, \dots, 0) \quad 1 \text{ at position } i$$

In the original space words are orthogonal.

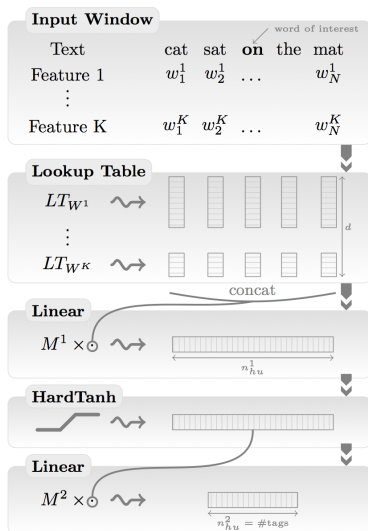
cat = (0,0,0,0,0,0,0,0,0,0,1,0,0,0,0, ...)

kitten = (0,0,1,0,0,0,0,0,0,0,0,0,0,0, ...)

To get the  $\mathbb{R}^d$  embedding vector for the word we multiply  $Wx$  where  $W$  is a  $d \times N$  vector with  $N$  words in the dictionary.



# Window Approach



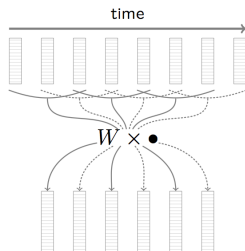
- Tags **one word** at the time
- Feed a **fixed-size window** of text around **each word** to tag
- **Works fine for most tasks**
- How do deal with **long-range dependencies**?

*E.g. in **SRL**, the **verb** of interest might be **outside** the **window**!*

## Sentence Approach

(1/2)

- Feed the **whole sentence** to the network
- Tag **one word** at the time: add extra **position** features
- **Convolutions** to handle variable-length inputs



See (Bottou, 1989)  
or (LeCun, 1989).

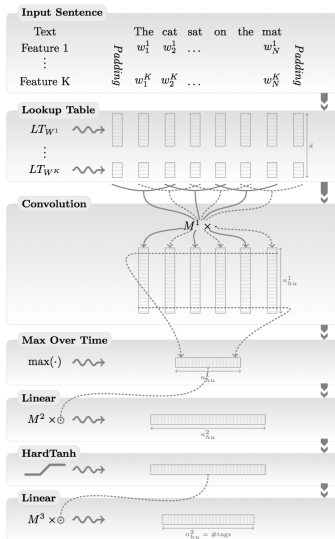
- Produces **local** features with higher level of abstraction
- **Max over time** to capture most relevant features



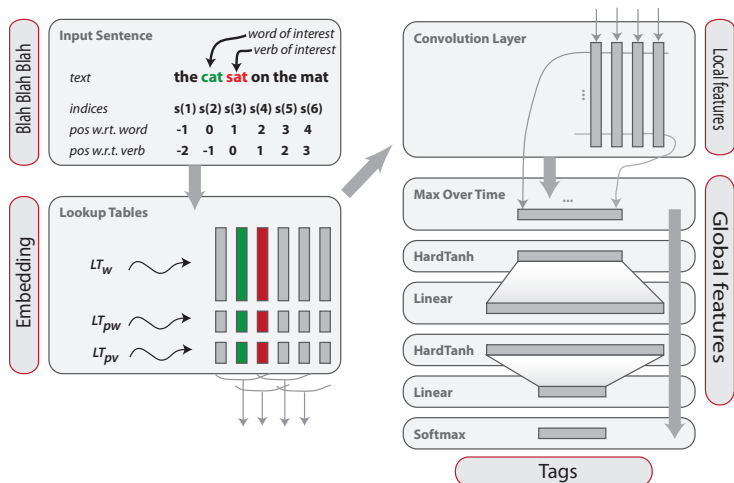
Outputs a **fixed-sized** feature vector

## Sentence Approach

(2/2)

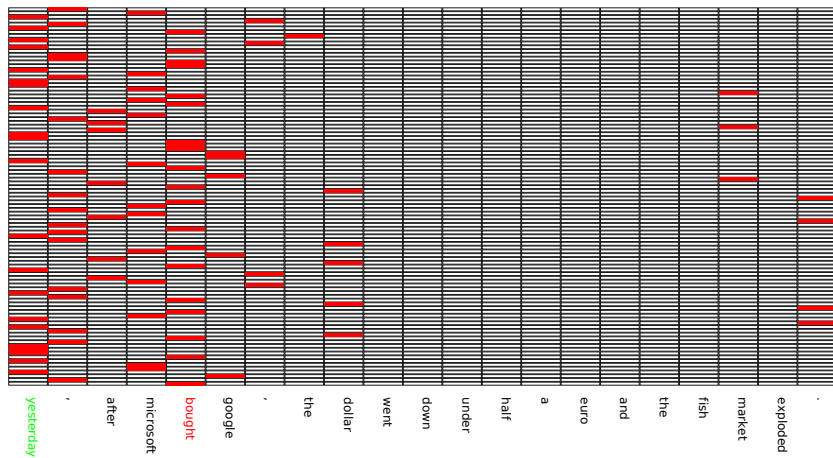


# Deep SRL

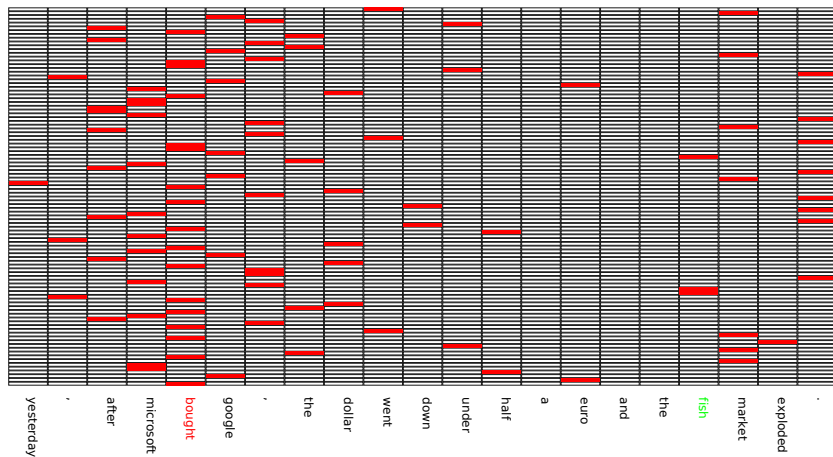


This is the network for a single window. We train/test predicting the entire sentence of tags ("structured outputs") using viterbi approach, similar to other NLP methods.

# Removing The Time Dimension (1/2)



# Removing The Time Dimension (2/2)



## Word Tag Likelihood (WTL)

- The network has one output  $f(\mathbf{x}, \mathbf{i}, \boldsymbol{\theta})$  per tag  $\mathbf{i}$
- Interpreted as a probability with a **softmax** over **all tags**

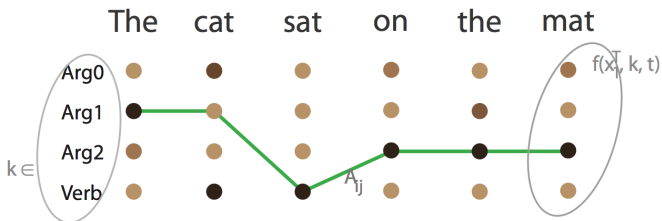
$$p(\mathbf{i} | \mathbf{x}, \boldsymbol{\theta}) = \frac{e^{f(\mathbf{x}, \mathbf{i}, \boldsymbol{\theta})}}{\sum_j e^{f(\mathbf{x}, \mathbf{j}, \boldsymbol{\theta})}}$$

*.. we can train directly for that (word tag likelihood) or we could train in a structured way by predicting the entire sentence's tags.*

That should be useful because tags are not independent.

## Sentence Tag Likelihood (STL)

- The **network score** for tag  $k$  at the  $t^{\text{th}}$  word is  $f([\mathbf{x}]_1^T, k, t, \boldsymbol{\theta})$
- $A_{kl}$  **transition score** to jump from tag  $k$  to tag  $l$



- Sentence** score for a tag path  $[i]_1^T$

$$s([\mathbf{x}]_1^T, [i]_1^T, \tilde{\boldsymbol{\theta}}) = \sum_{t=1}^T \left( A_{[i]_{t-1}[i]_t} + f([\mathbf{x}]_1^T, [i]_t, t, \boldsymbol{\theta}) \right)$$



## Supervised Benchmark Results

- Network architectures:
  - ★ Window (5) approach for POS, CHUNK & NER (300HU)
  - ★ Convolutional (3) for SRL (300+500HU)
  - ★ Word Tag Likelihood (WTL) and Sentence Tag Likelihood (STL)
- Network features: lower case words (size 50), capital letters (size 5)  
dictionary size 100,000 words

Approach	POS (PWA)	Chunking (F1)	NER (F1)	SRL (F1)
<b>Benchmark Systems</b>	<b>97.24</b>	<b>94.29</b>	<b>89.31</b>	<b>77.92</b>
NN+WTL	96.31	89.13	79.53	55.40
NN+STL	96.37	90.33	81.47	70.99

- STL helps, but... fair performance.
- Capacity mainly in words features... are we training it right?

## Supervised Word Embeddings

- Sentences with **similar words** should be **tagged in the same way**:

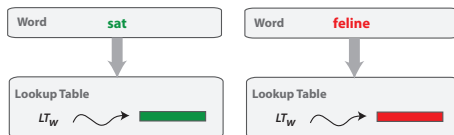
- ★ The **cat** sat on the mat
- ★ The **feline** sat on the mat

france	jesus	xbox	reddish	scratched	megabits
454	1973	6909	11724	29869	87025
persuade	thickets	decadent	widescreen	odd	ppa
faw	savary	divo	antica	anchieta	uddin
blackstock	sympathetic	verus	shabby	emigration	biologically
giorgi	jfk	oxide	awe	marking	kayak
shaheed	khwarazm	urbina	thud	heuer	mclarens
rumelia	stationery	epos	occupant	sambhaji	gladwin
planum	ilias	eglinton	revised	worshippers	centrally
goa'uld	gsNUMBER	edging	leavened	ritsuko	indonesia
collation	operator	frg	pandionidae	lifeless	moneo
bacha	w.j.	namsos	shirt	mahan	nilgiris

- About **1M** of words in WSJ
- **15%** of **most frequent words** in the dictionary are seen **90%** of the time
- **Cannot expect words to be trained properly!**

# Improving Word Embedding

- Rare words are not trained properly
- Sentences with similar words should be tagged in the same way:
  - The cat sat on the mat
  - The feline sat on the mat



Only 1M WSJ not enough – let's use lots of unsupervised data!

## Semi-supervised: MTL with Unlabeled Text

- Language Model: “*is a sentence actually english or not?*”  
Implicitly captures:      \* syntax      \* semantics
- Bengio & Ducharme (2001) Probability of next word given previous words. Overcomplicated – we do not need probabilities here
- English sentence windows: Wikipedia ( $\sim 631M$  words)  
Non-english sentence windows: middle word randomly replaced  
     the champion federer wins wimbledon  
     vs. the champion saucepan wins wimbledon
- Multi-class margin cost:

$$\sum_{s \in \mathcal{S}} \sum_{w \in \mathcal{D}} \max(0, 1 - f(s, w_s^*) + f(s, w))$$

$\mathcal{S}$ : sentence windows     $\mathcal{D}$ : dictionary

$w_s^*$ : true middle word in  $s$

$f(s, w)$ : network score for sentence  $s$  and middle word  $w$

# Language Model: Embedding

Nearest neighbors in 100-dim. embedding space:

FRANCE 454	JESUS 1973	XBOX 6909	REDDISH 11724	SCRATCHED 29869
SPAIN	CHRIST	PLAYSTATION	YELLOWISH	SMASHED
ITALY	GOD	DREAMCAST	GREENISH	RIPPED
RUSSIA	RESURRECTION	PSNUMBER	BROWNISH	BRUSHED
POLAND	PRAYER	SNES	BLUISH	HURLED
ENGLAND	YAHWEH	WII	CREAMY	GRABBED
DENMARK	JOSEPHUS	NES	WHITISH	TOSSED
GERMANY	MOSES	NINTENDO	BLACKISH	SQUEEZED
PORTUGAL	SIN	GAMECUBE	SILVERY	BLASTED
SWEDEN	HEAVEN	PSP	GREYISH	TANGLED
AUSTRIA	SALVATION	AMIGA	PALER	SLASHED

(Even fairly rare words are embedded well.)

# Results

Algorithm	POS (PWA)	CHUNK (F1)	NER (F1)	SRL (F1)
Baselines	97.24 [Toutanova '03]	94.29 [Sha '03]	89.31 [Ando '05]	77.92 [Koomen '05]
NN + WTL	96.31	89.13	79.53	55.40
NN + STL	96.37	90.33	81.47	70.99
NN + LM + STL	97.22	94.10	88.67	74.15
NN + ... + tricks	97.29 [+suffix]	94.32 [+POS]	89.95 [+gazetteer]	76.03 [+Parse Trees]

## NOTES:

- Didn't compare to benchmarks that used external labeled data.
- [Ando '05] uses external unlabeled data.
- [Koomen '05] uses 4 parse trees not provided by the challenge. Using only 1 tree it gets 74.76.

# Software

Code for tagging with POS, NER, CHUNK, SRL + parse trees:

<http://ml.nec-labs.com/senna/>

System	RAM (Mb)	Time (s)
Toutanova, 2003	1100	1065
Shen, 2007	2200	833
SENNA	32	4

(a) POS

System	RAM (Mb)	Time (s)
Koomen, 2005	3400	6253
SENNA	124	52

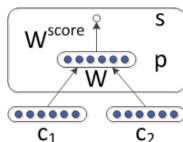
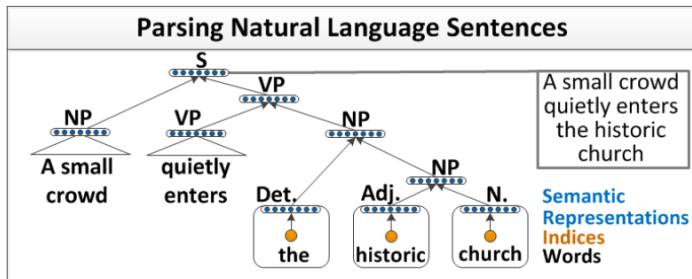
(b) SRL

See also Torch: <http://www.torch.ch>

# Recursive NNs for Parsing, Sentiment, ... and more!

(Socher et al., ICML '13), (Socher et al., EMNLP, '13))

Build sentence representations using the parse tree to compose embeddings via a nonlinear function taking pairs  $(c_1, c_2)$  and output  $p$ .



$$s = W^{score} p \quad (9)$$

$$p = f(W[c_1; c_2] + b)$$

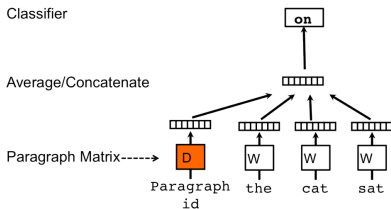
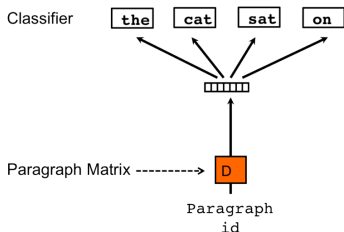


# Paragraph Vector

(Le & Mikolov, '14)

*A Paragraph Vector (a vector that represents a paragraph/doc) learned by:*

- 1) Predicting the words in a doc;
- 2) predict  $n$ -grams in the doc:



At test time, for a new document, one needs to learn its vector, this can encode word order via the  $n$ -gram prediction approach.

# Comparison of CNN, RNN & PV (Kim '14)

Model	MR	SST-1	SST-2	Subj	TREC	CR	MPQA
CNN-rand	76.1	45.0	82.7	89.6	91.2	79.8	83.4
CNN-static	81.0	45.5	86.8	93.0	92.8	84.7	<b>89.6</b>
CNN-non-static	<b>81.5</b>	48.0	87.2	93.4	93.6	84.3	89.5
CNN-multichannel	81.1	47.4	<b>88.1</b>	93.2	92.2	<b>85.0</b>	89.4
RAE (Socher et al., 2011)	77.7	43.2	82.4	—	—	—	86.4
MV-RNN (Socher et al., 2012)	79.0	44.4	82.9	—	—	—	—
RNTN (Socher et al., 2013)	—	45.7	85.4	—	—	—	—
DCNN (Kalchbrenner et al., 2014)	—	48.5	86.8	—	93.0	—	—
Paragraph-Vec (Le and Mikolov, 2014)	—	<b>48.7</b>	87.8	—	—	—	—
CCAE (Hermann and Blunsom, 2013)	77.8	—	—	—	—	—	87.2
Sent-Parser (Dong et al., 2014)	79.5	—	—	—	—	—	86.3
NBSVM (Wang and Manning, 2012)	79.4	—	—	93.2	—	81.8	86.3
MNB (Wang and Manning, 2012)	79.0	—	—	<b>93.6</b>	—	80.0	86.3
G-Dropout (Wang and Manning, 2013)	79.0	—	—	93.4	—	82.1	86.1
F-Dropout (Wang and Manning, 2013)	79.1	—	—	<b>93.6</b>	—	81.9	86.3
Tree-CRF (Nakagawa et al., 2010)	77.3	—	—	—	—	81.4	86.1
CRF-PR (Yang and Cardie, 2014)	—	—	—	—	—	82.7	—
SVM <sub>S</sub> (Silva et al., 2011)	—	—	—	—	<b>95.0</b>	—	—

Table 2: Results of our CNN models against other methods. **RAE**: Recursive Autoencoders with pre-trained word vectors from Wikipedia (Socher et al., 2011). **MV-RNN**: Matrix-Vector Recursive Neural Network with parse trees (Socher et al., 2012). **RNTN**: Recursive Neural Tensor Network with tensor-based feature function and parse trees (Socher et al., 2013). **DCNN**: Dynamic Convolutional Neural Network with k-max pooling (Kalchbrenner et al., 2014). **Paragraph-Vec**: Logistic regression on top of [paragraph](#) vectors (Le and Mikolov, 2014). **CCAE**: Combinatorial Category Autoencoders with combinatorial category grammar operators (Hermann and Blunsom, 2013). **Sent-Parser**: Sentiment analysis-specific parser (Dong et al., 2014). **NBSVM**, **MNB**: Naive Bayes SVM and Multinomial Naive Bayes with uni-bigrams from Wang and Manning (2012). **G-Dropout**, **F-Dropout**: Gaussian Dropout and Fast Dropout from Wang and Manning (2013). **Tree-CRF**: Dependency tree

## Some More Recent Work

- Compositionality approaches by Marco Baroni's group:  
Words are combined with linear matrices dependent on the P.O.S.:  
G. Dinu and M. Baroni. How to make words with vectors: Phrase generation in distributional semantics. ACL '14.
- Document representation by Phil Blunson's group:  
Variants of convolutional networks for text:  
Kalchbrenner et al. A Convolutional Neural Network for Modelling Sentences. ACL '14

Good tutorial slides from these teams covering multiple topics:

New Directions in Vector Space Models of Meaning

<http://www.cs.ox.ac.uk/files/6605/aclVectorTutorial.pdf>

# Summary

- Generic end-to-end deep learning system for NLP tasks
- Word embeddings combined to form sentence or document embeddings can perform well on supervised tasks.
- Previous common belief in NLP: engineering syntactic features necessary for semantic tasks.  
One can do well by engineering a model/algorithm rather than features.

Attitude is changing in recent years... let's see what happens!